

Integrating Node-Place Model With Shapley Additive Explanation for Metro Ridership Regression

Yizhe Wang^{1b}, Zijia Wang^{1b}, and Fanxi Zhao^{1b}

Abstract—Accurate metro ridership estimation is essential for effective urban transportation planning. Traditional regression methods relying on historical passenger flow data often overlook the intrinsic relationship between network scale and land use, limiting predictive accuracy. To address this issue, this study focuses on the Beijing urban rail transit system, aiming to explore how station-level factors influence spatial variations in metro ridership over a decade (2013–2022). An advanced machine learning framework that integrates Node-Place model with Shapley Additive Explanation (SHAP-based machine learning method) is proposed. Twelve influencing factors are developed, including network centrality measures and points of interest (POI) attributes, and a CRITIC weighting method is applied to weight these variables within the Node-Place Model, ensuring an objective assessment of their relative importance. Utilizing LightGBM regression method combined with SHAP values, our model achieved an average R^2 of 0.88 in annual average daily entry ridership regression and a smaller discrepancy between R^2 and adjusted R^2 , significantly outperforming ordinary least squares (OLS, $R^2 = 0.42$) and geographically weighted regression (GWR, $R^2 = 0.58$). The proposed model consistently demonstrates an AIC value approximately 15% lower than alternative models across multiple years of regression tasks, highlighting its stability and superior performance. Key indicators identified include PageRank and the number of restaurants and enterprises around metro stations. Combining detailed land use data and network centrality measures with advanced machine learning techniques proves to be an effective way to enhancing ridership regression.

Index Terms—Urban rail transit, node-place model, critic weighting method, SHAP values.

I. INTRODUCTION

GLOBALLY, mounting concerns over traffic congestion, environmental degradation, and urban sustainability have encouraged governments and transportation authorities to adopt more holistic and integrated approaches to urban mobility [1], [2], [3]. The conventional over-reliance on private vehicles not only imposes economic and environmental burdens, but also hinders the creation of livable, well-functioning

cities. Against this backdrop, urban rail transit (URT) systems have emerged as a preferred mobility solution, demonstrating remarkable potential in mitigating congestion, curbing emissions, and shaping sustainable urban growth patterns.

In China, the development of URT networks has been particularly vigorous, with more than 30 major cities—including Shanghai, Shenzhen, and Wuhan—demonstrating a rich history and extensive experience in planning, constructing, and expanding rail infrastructures [4], [5], [6]. These efforts have substantially alleviated urban traffic problems and offered insightful examples of integrated transport and land use strategies. A prominent instance is the Beijing metro system, which has witnessed rapid growth over the past decade. By April 2023, the Beijing network extended beyond 800 kilometers and encompassed 27 lines, establishing itself as an expansive transit system. This remarkable expansion underscores China's dedication to strategic urban transport planning and infrastructure development, while also plays a significant role in influencing international trends in the evolution of urban rail transit systems [7].

Network expansion not only enhances accessibility but also acts as a catalyst for the development of surrounding land areas, exerting a substantial influence on passenger flow at existing stations [8], [9]. Traditional studies frequently attribute variations in public transport ridership to factors such as economic conditions, urban development dynamics, population shifts, and changes in fare structures and service offerings [10]. However, in the context of the rapid urbanization characteristic of Chinese cities, comprehensive and up-to-date data on land use and other related variables are often scarce or outdated. Consequently, prevailing research methodologies predominantly rely on historical passenger flow data to analyze ridership changes, thereby inadequately capturing the intricate interplay between network topology and land use variations. This limitation highlights the necessity for more sophisticated analytical frameworks that integrate detailed land use information and network structural attributes to more accurately understand and predict ridership patterns.

To address this issue, we employ a combination of the Node-Place model [11] and a SHAP-based machine learning interpretability model [12] to quantify the influence of various indices on dynamic changes in metro ridership. For the selection of measurement indicators and characteristics of stations within the entire metro network, this study accounts

Received 11 December 2024; revised 12 February 2025; accepted 21 February 2025. This work was supported by the Beijing Natural Science Foundation under Grant L221025. The Associate Editor for this article was S. Ahmad. (Corresponding author: Zijia Wang.)

Yizhe Wang and Zijia Wang are with the School of Civil Engineering, Beijing Jiaotong University, Beijing 100044, China (e-mail: 24110349@bjtu.edu.cn; zjwang@bjtu.edu.cn).

Fanxi Zhao is with the Department of Statistics and Actuarial Science, Faculty of Mathematics, University of Waterloo, Waterloo, ON N2L 3G1, Canada (e-mail: f52zhao@uwaterloo.ca).

Digital Object Identifier 10.1109/TITS.2025.3546471

for the features of the surrounding built environment, including land use patterns and housing prices. Drawing upon a comprehensive review of historical literature, we develop twelve explanatory variables that potentially influence the annual average daily entry passenger flow at metro stations.

In this study, we focus on the Beijing urban rail transit system, aiming to elucidate how station-level factors drive spatial variations in metro ridership. We introduce an advanced machine learning framework that integrates traditional transportation models with SHAP-based techniques, enhancing both the accuracy and interpretability of ridership estimation and regression. By incorporating an objective weighting method, we quantitatively assign weights to key indicators, ensuring a more precise evaluation of station characteristics. Additionally, critical factors influencing ridership patterns are identified, providing valuable insights into the complex dynamics of urban rail transit.

II. LITERATURE REVIEW

In the study of the synergistic evolution and interplay between transit networks and passenger flow, one notable challenge is the difficulty in accessing comprehensive land use type data [13]. Addressing this issue, Enjian Yao et al. innovate a proxy variable approach for land use categorization. This method involves using key indicators such as Morning Peak Hour Boarding Factor (MPHFB), Morning Peak Hour Alighting Factor (MPHFA), Evening Peak Hour Boarding Factor (EPHFB), and Evening Peak Hour Alighting Factor (EPHFA) for clustering analysis [14]. The objective is to discern patterns in the land use characteristics surrounding transit stations with similar passenger flow profiles. This approach significantly alleviates the extensive effort required in gathering detailed land use data. However, given their reliance on area-based metrics, it's crucial to acknowledge that inferences based on this proxy, particularly those concerning spatial congruence, may not align with the on-ground realities.

Building upon the challenges of integrating land use data with transit network analysis, the Node-Place (NP) Model [15] offers a robust framework for evaluating the relationship between metro network nodes and Points of Interest (POI) values. Initially proposed by Bertolini, the Node-Place Model serves as a regionally oriented planning approach that emphasizes the interplay between transportation infrastructure and land use, guided by public transit principles [16]. This model conceptualizes Transit-Oriented Development (TOD) coverage areas as both nodes within the transportation network and integral components of urban space [17]. By doing so, the Node-Place Model advocates for enhancing regional transportation support to promote both the density and diversity of land use, thereby fostering a mutually reinforcing dynamic between transit infrastructure and urban development.

The Node-Place Model not only provides a theoretical foundation for assessing the temporal and qualitative aspects of transportation network construction but also offers practical insights for TOD-oriented planning. By emphasizing the symbiotic relationship between transit nodes and their surrounding urban areas, the Node-Place Model facilitates a comprehensive evaluation of transportation and urban development levels,

their degree of matching, and interactive potential. This dual perspective allows for a micro-level exploration of the coordination between transit functions and land use within station regions. Consequently, the Node-Place Model supports the creation of dynamic systems where transportation enhancements lead to increased land use density and diversity, promoting sustainable urban growth [18].

Furthermore, empirical studies leveraging the Node-Place Model have demonstrated its efficacy in various urban contexts. For instance, Guowei Lyu et al. utilized the model to develop a TOD typology for Beijing metro station areas, highlighting its adaptability and relevance in rapidly urbanizing environments. These applications underscore the model's capacity to guide urban planners in making informed decisions that balance transportation efficiency with urban livability [19].

In contemporary machine learning applications, comprehending the rationale behind a model's predictions is as crucial as the accuracy of those predictions. This is particularly evident in tasks involving complex datasets, where increasing model complexity often leads to significant enhancements in predictive performance. However, such advancements frequently result in models that are challenging for experts to interpret. For instance, deep learning and ensemble learning models, despite their superior predictive capabilities, exhibit "black-box" characteristics that obscure their decision-making processes [20].

Despite the proliferation of explainability methods, many of these approaches suffer from a lack of clear interrelations and comprehensive frameworks that guide practitioners on their optimal application contexts.

To address these challenges, SHAP (SHapley Additive exPlanations) has emerged as a unified framework that effectively interprets machine learning model predictions by leveraging concepts from cooperative game theory, specifically Shapley values. SHAP assigns each feature an importance value that quantifies its contribution to the prediction, offering several notable advantages over other feature importance methods. These include consistency—ensuring that if a feature's contribution increases, its SHAP value does not decrease—and fairness, where the sum of SHAP values equals the model's baseline output in the absence of any feature information. Additionally, SHAP provides both local explanations for individual predictions and global explanations that reveal overall feature importance trends across the entire dataset [21].

The robustness and versatility of SHAP have made it a preferred choice for feature importance measurement across various domains. In the transportation domain, the application of SHAP is gaining traction, albeit still in its early stages. For instance, Oseni et al. utilized SHAP within a deep learning framework for resilient intrusion detection in IoT-enabled transportation networks, successfully identifying the most influential features contributing to security threats [20]. Similarly, Li et al. introduced SVCE (Shapley Value Guided Counterfactual Explanation) for autonomous driving systems, which leverages SHAP to generate counterfactual explanations, thereby enhancing model interpretability and reliability [22]. Besides, SHAP has been demonstrated to enhance the accuracy of passenger flow predictions, optimize

TABLE I
OVERVIEW OF THE LIMITATIONS IN RELATED STUDIES

Source	Methods	Key Findings	Limitations
Enjian Yao et al. [14]	AFC data Driven	Explored the heterogeneity of the impact of metro network expansion and land use on ridership growth.	Proxy-based inference may compromise prediction accuracy.
Bertolini et al. [16][17][18]	Traditional Node-Place Model	Developed a foundational model for transportation supply-demand balance analysis.	Qualitative analysis with subjective influence, lacking support from quantitative analysis.
	OLS(ordinary least squares)	Identified linear relationships(positive/negative) between built environment factors and metro ridership.	Lacks consideration of spatial data features, leading to the loss of transportation insights.
Mengya Li et al. [25]	GWR(geographically weighted regression)	Revealed spatial heterogeneity in the impact of various factors on metro ridership.	Sensitivity to the selection of spatial weights and limited interpretability.

resource allocation, and improve energy management in urban rail transit, contributing to reduced operational costs, increased efficiency, and enhanced sustainability [23]. By identifying the key factors influencing travel patterns, SHAP facilitates dynamic scheduling and more efficient capacity management. In the context of intelligent driving, SHAP has proven effective in improving pedestrian intent prediction and enhancing the transparency of decision-making processes, thereby increasing safety, reducing accidents, and lowering associated costs [24]. Furthermore, SHAP contributes to the optimization of energy consumption and driving strategies, leading to additional reductions in operational expenses. These studies demonstrate SHAP's potential to provide actionable insights in critical areas of intelligent transportation systems (ITS), where transparency and reliability are paramount.

Building upon previous research, we assert that employing regression methods for the accurate and reasonable estimation of metro ridership at metro stations with specific land use characteristics can yield significant economic benefits [25]. A summary of relevant studies, along with their limitations, is provided in Table I.

III. METHODOLOGY AND INDEX DEFINITION

In order to evaluate the evolution of factors influencing station development and the resultant spatial variations in metro ridership, we adopt a multi-faceted methodological approach that integrates the Node-Place Model with SHAP-based machine learning techniques. The flowchart of the proposed method is shown in Fig. 1.

The process begins by incorporating network data and POI characteristics as inputs. The Node-Place Model is then applied to assess both node and place indices, with weights objectively assigned through the CRITIC weighting method [26]. A machine learning model is subsequently employed to perform metro ridership regression, with SHAP values offering interpretability by quantifying the contribution of each feature. Finally, the model's performance is evaluated across multiple regression tasks by proposed metrics.

A. Data Description And Index Definition

For this study, we collected annual average daily entry passenger flow data for all Beijing Metro stations

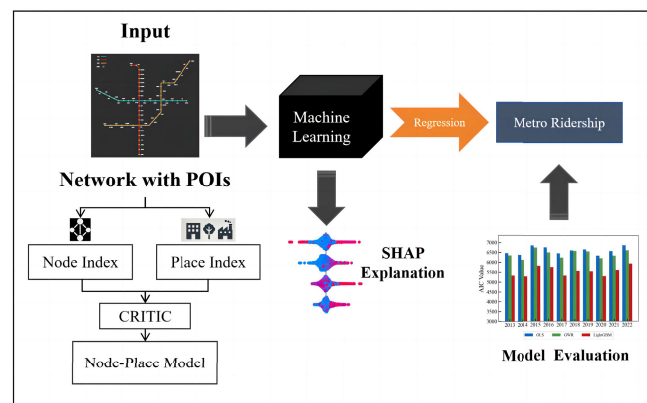


Fig. 1. Flowchart of proposed method.

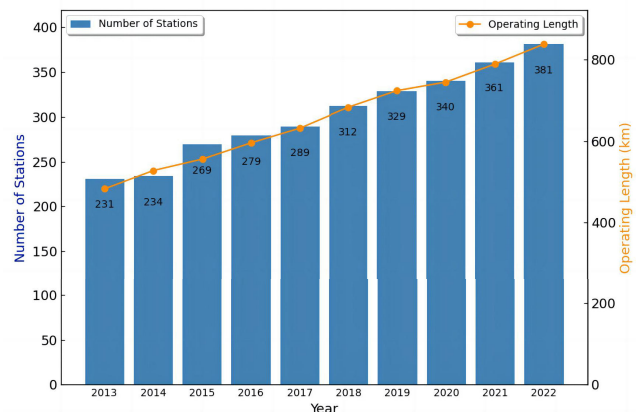


Fig. 2. Beijing metro network expansion.

from 2013 to 2022, which serves as the dependent variable in our analysis. The number of stations included in our dataset is consistent with the station counts illustrated in Fig. 2.

Complex network theory provides a new perspective for characterizing the topological features of transportation networks. Traditionally, degree centrality, closeness centrality, and betweenness centrality are commonly used to represent the connectivity, accessibility, and transitivity of a city within its transportation network, respectively [27]. In this study, we build upon these centrality measures and substitute degree centrality with the adapted PageRank algorithm to address the

TABLE II
SUMMARY OF VARIABLE DESCRIPTIONS

Variable	Description	Shorthand notation	Feature Type
PageRank	Network node importance measure based on the adapted PageRank algorithm	N1	Topological attributes
Betweenness	Network node centrality measure based on the counts of the shortest paths that traverse each station	N2	Topological attributes
Closeness	Network node accessibility measure based on the speed at which a station can connect to every other station	N3	Topological attributes
Eigenvector	Network node importance measure based on the quantity and quality of its neighboring nodes	N4	Topological attributes
Shopping	Business activity level, Number of shopping POIs within 800-meter	P1	Numerical
Enterprise	Enterprise potential, Number of enterprise POIs within 800-meter	P2	Numerical
Residence	Residential density, Number of residence POIs within 800-meter	P3	Numerical
Accommodation	Accommodation availability, Number of accommodation POIs within 800-meter	P4	Numerical
Restaurant	Dining activity level, Number of restaurant POIs within 800-meter	P5	Numerical
Hospital	Healthcare accessibility, Number of hospital POIs within 800 meters	P6	Numerical
Sale	Average house price within 800-meter (CNY/m ²)	P7	Price(Numerical)
Rent	Average rent price within 800-meter (CNY/m ² /month)	P8	Price(Numerical)

issue of homogeneity in degree centrality. PageRank allows for a differentiated representation of node importance, accounting for both the quantity and quality of connections. Additionally, we incorporate eigenvector centrality, which simultaneously captures the importance of a node and its neighbor [28]. These four topological indices together form a comprehensive evaluation framework that can provide a multidimensional representation of the node characteristics of stations within the network.

Regarding the land use data, an 800-meter (10-minute walking) buffer zone has been validated as an effective catchment area for metro stations [25]. Consequently, the number of POIs from various categories is calculated within this 800-meter radius around each station, which captures the station's location significance within the Node-Place Model framework, providing a nuanced representation of its spatial and functional attributes.

Finally, in alignment with the principles of scientific rigor, practical applicability, and completeness, Table II presents a carefully selected set of 12 indicators for the Node-Place model, encompassing four classic node-based centrality metrics and eight POI-related features based on history research [29], [30]. These indicators capture both the topological characteristics of the stations and the surrounding land use patterns, offering a robust framework for evaluating the performance of the metro network.

B. Node Place Model

Within the Node-Place model, the synergistic effects between nodes (stations) and places (areas of interest) are qualitatively assessed using a Cartesian coordinate system. As illustrated in Fig. 3, the distribution of the Node Place

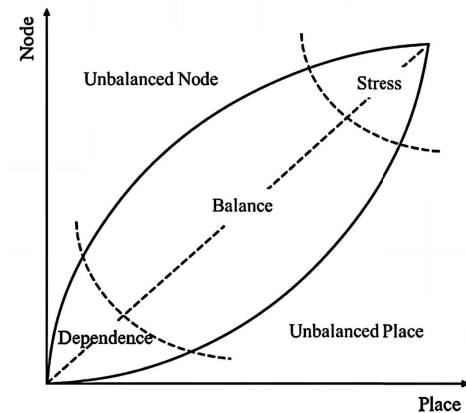


Fig. 3. Node-place model proposed by bertolini.

model along the central dashed line represents an ideal balance between transportation performance and land use efficiency. The spindle-shaped regions of the model are categorized into three distinct states—Dependence, Balance, and Stress—based on the combined levels of node value and place value.

Specifically, Stress state located at the top of the spindle, indicating scenarios where both transportation attractiveness and land development intensity are high. In this state, the interaction potential approaches saturation, and further growth may lead to incompatibilities between transportation and urban development. Dependence State positioned at the bottom of the spindle, signifying areas where the potential for development remains underexploited. This suggests that the region has untapped growth opportunities that have yet to be fully realized. Balance State situated between the Stress and Dependence states, reflecting a relative equilibrium between

node value and place value. In this state, the value of the node and the place are harmoniously aligned, indicating a well-integrated transportation and land use system.

Additionally, the lower right quadrant of the spindle represents the Unbalanced Place state, where urban development has outpaced transportation services (e.g., CBD areas). In this scenario, the place value significantly exceeds the node value, highlighting a disparity between land use and transportation infrastructure, with transportation facilities falling short of demand. Conversely, the upper left quadrant denotes the Unbalanced Node state, where transportation services surpass the demands of urban development (e.g., peripheral station areas). Here, the node value is markedly higher than the place value, indicating that the development of transportation infrastructure exceeds urban transportation demand, thereby demonstrating a surplus in supply.

C. CRITIC Weighting Method

To effectively incorporate multiple indicators within the Node-Place framework, we employ the CRITIC weighting method. This integration transforms the Node-Place model from a purely qualitative approach into a quantitative one by objectively assigning weights to each indicator based on their inherent attributes.

The CRITIC weighting method employs the Pearson correlation coefficient to evaluate the relationships among indices, thereby determining the degree of conflict between indicators, which is defined as follows:

$$r_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}} \quad (1)$$

where $\text{Cov}(X, Y)$ denotes the covariance between variables X and Y , $\text{Var}(X)$ represents the variance of X , and $\text{Var}(Y)$ signifies the variance of Y . Finally, the correlation coefficient is employed to quantify the conflict between indicators:

$$R_j = \sum_{i=1}^n (1 - r_{ij}), \quad j = 1, 2, \dots, n \quad (2)$$

where r_{ij} denotes the Pearson correlation coefficient that quantifies the relationship between evaluation index i and index j .

Based on the derived index conflict, the information content of each indicator can subsequently be calculated:

$$C_j = \sigma_j \times R_j, \quad j = 1, 2, \dots, n \quad (3)$$

where σ_j denotes the standard deviation of the evaluation index j . Ultimately, C_j quantifies the significance of the j -th evaluation index within the entire evaluation system. We employ above formula to calculate the initial weight of each indicator within its respective category (e.g., Betweenness within the Node index and Shopping within the Place index). Finally, the objective weight assigned to the j -th evaluation index is defined by:

$$W_j = \frac{C_j}{\sum_{j=1}^n C_j} \quad (4)$$

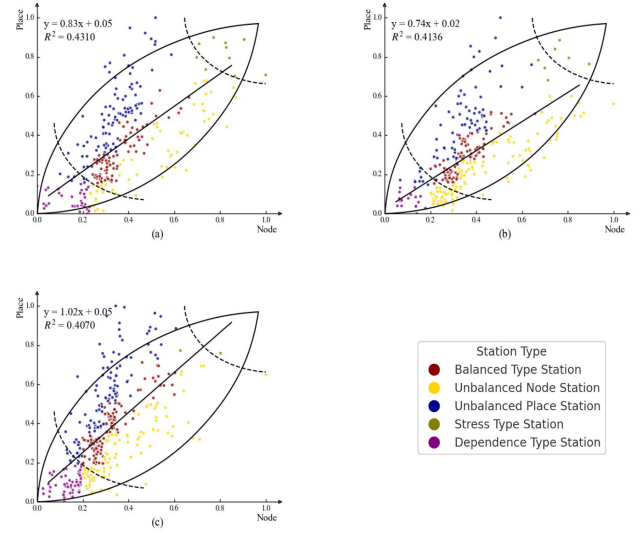


Fig. 4. (a)Node-place indexes for 2016; (b)Node-place indexes for 2019; (c)Node-place indexes for 2022.

D. SHapley Additive Explanation

The SHAP model is a method based on game theory and local explanations, designed to estimate the marginal contribution of different features to the output prediction. This method involves the computation of Shapley values, which represent the contributions of the features as follows.

$$\phi_k = \sum_{z' \subseteq x'} \frac{|z'|!(K - |z'| - 1)!}{K!} \left[f_x(z') - f_x\left(\frac{z'}{k}\right) \right] \quad (5)$$

In (5), ϕ_k represents the contribution value of the k -th feature, indicating its influence on the model's prediction. K denotes the total number of features included in the model, encompassing all relevant indicators under consideration. x refers to the feature data point, encapsulating the set of input variables used for generating predictions. $|z'|$ signifies the number of non-zero terms within the feature vector. $f(x)$ denotes the predicted output of the model for the given input data point, representing the model's forecast based on the provided features. Specifically, in this paper, $K = 12$ represents the total number of features, and x corresponds to the set of input variables (e.g., network centrality and POI attributes), with $f(x)$ denoting the predicted ridership output based on input features.

IV. RESULTS AND ANALYSIS

A. CRITIC - Node-Place Model - R^2

Employing the CRITIC weighting method, we allocate weights to the indicators listed in TABLE II. The weighted outcomes are subsequently normalized for visualization through the Node-Place model. To evaluate the evolution of factors influencing station development, we employed Ordinary Least Squares regression to conduct a longitudinal analysis of the Node-Place model for Beijing rail transit system over a ten-year period, analyzing the balance between node values and place values over years [18]. As depicted in Fig. 4, the model consistently achieved an R-squared value of approximately 0.41, indicating a stable and reliable capacity to explain the variance.

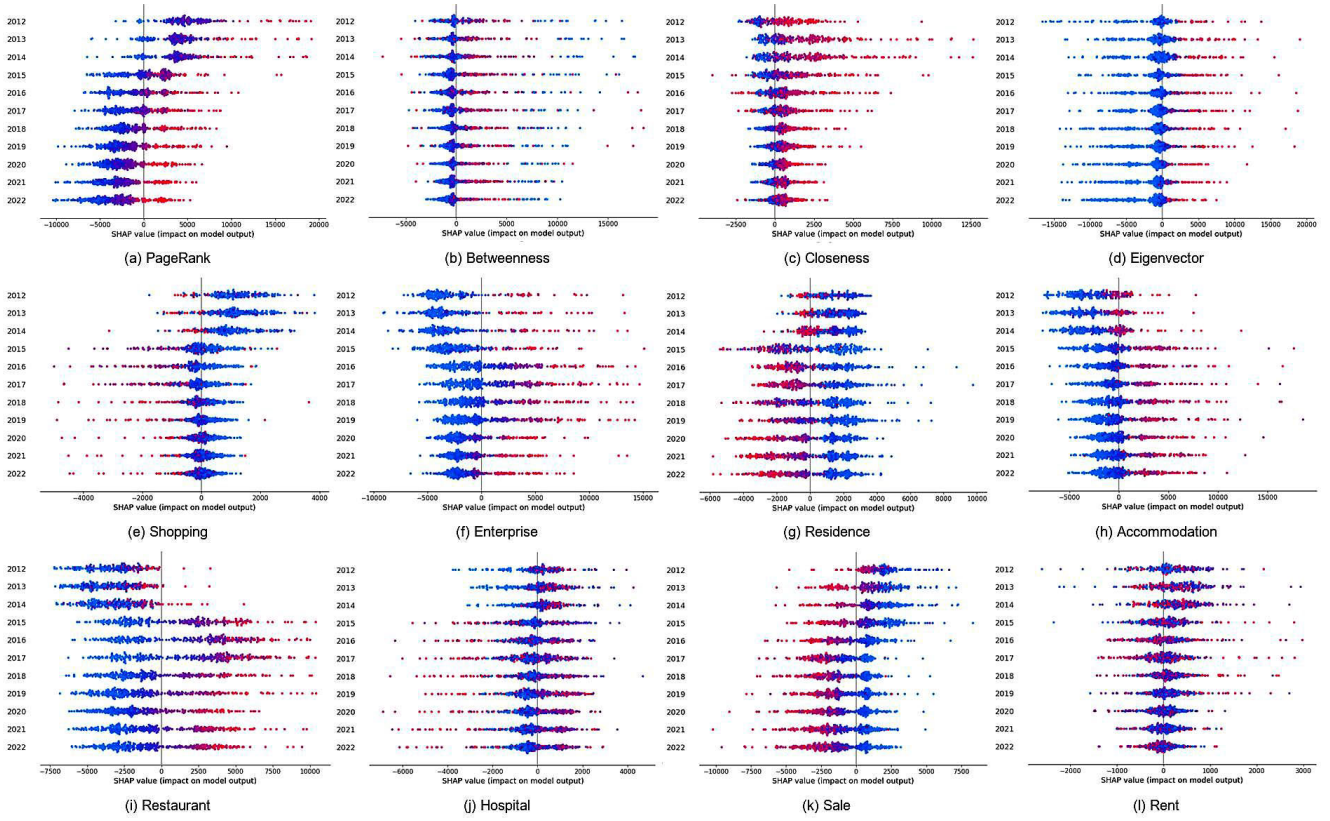


Fig. 5. SHAP values' beeswarm plot of the total 12 node-place features.

From the perspective of slope coefficients, the year 2016 represented the early stages of network development, with the metro system's scale not yet fully realized. The recorded slope value of 0.83, which is below the unit benchmark of 1, suggests that the network was still in its nascent phase, characterized by node-dominated features. As the network expansion accelerated, the slope in 2019 decreased further to 0.74, indicating that the growth of Place Value was lagging behind the rapid development of the network. By 2022, this trend appeared to reverse, with the slope increasing to 1.02, signifying a rise in Place Value and indicating a shift towards a transportation-oriented development model. Throughout these three years, the model's y-intercept remained consistently close to zero, reinforcing the model's reliability in capturing the dynamic interplay between network expansion and place-based factors.

B. Explanations Based on SHAP Value

To elucidate the spatial variations in metro ridership attributable to station-level factors, existing research predominantly utilizes Ordinary Least Squares or Geographically Weighted Regression (GWR) models to examine the influence of built environment elements on rail transit passenger flows, these conventional methods often fall short in adequately capturing the spatio-temporal heterogeneity inherent in the data and exhibit limited explanatory power [25].

In response to these limitations, we integrate a multiple regression method based on LightGBM [31], [32], a highly efficient gradient boosting decision tree algorithm, with Shapley additive explanation based on SHAP value to provide a

more comprehensive analysis of how Node-Place indicators within Beijing's rail transit stations contribute to the spatial dynamics of ridership.

To investigate the relative changes in feature importance over the past decade, we computed the SHAP values for node and place indicators annually using (5), with the results depicted in Fig. 5. The SHAP values are visualized as 'contribution amounts', where the distance between each station's scatter point and the zero baseline reflects the magnitude of the SHAP value. Positive SHAP values, located to the right of the baseline, indicate a positive impact on ridership predictions, while negative values to the left signify a negative influence. Additionally, the color of each scatter point represents the original value of the input feature, with red indicating higher values and blue indicating lower values, and gradient colors representing intermediate values. Consequently, the contributions of the twelve NP indicators for each station may offset each other in the prediction output.

Fig. 5 reveals significant variations in the contributions of different NP indicators across stations and over time. Notably, indicators such as 'Hospital' and 'Eigenvector' exhibit substantial differences in their influence on ridership across different stations within the same period. Furthermore, the long-term impact of certain indicators, like 'PageRank', has shifted from positive to negative over the past decade.

During the network's overall evolution, some indicators demonstrated considerable changes in their influence on ridership; for instance, the 'Restaurant' indicator initially had a significant negative contribution to ridership predictions across

TABLE III
EVOLUTION OF VARIABLES' SHAP VALUE

Variables	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022
N1	4490.31	4337.99	2390.22	2554.51	2430.71	2594.13	2947.26	2998.10	3317.39	3627.68
N2	2011.19	2026.61	1905.44	1813.84	1724.95	1650.85	1545.36	1373.02	1539.23	1361.94
N3	2112.90	2156.01	1375.91	1154.93	886.87	756.57	738.65	555.14	572.53	585.00
N4	2815.12	2809.55	2505.59	2437.55	2385.18	2672.42	2530.75	2478.24	2546.21	2841.71
P1	1236.37	1030.13	683.11	738.75	692.94	497.60	477.45	373.60	372.50	367.72
P2	3870.12	3868.61	3086.41	3088.96	3012.93	2644.69	2798.10	2271.84	2170.08	2138.11
P3	1420.14	1295.50	1905.44	1623.03	1697.63	1749.09	1734.13	1700.90	1798.92	1703.74
P4	3018.59	2733.87	2570.01	2071.85	2136.28	2342.84	2199.84	1906.55	2015.88	1905.32
P5	3466.22	3105.00	2954.49	3461.22	3448.18	2748.14	2891.97	2557.90	2751.01	2582.01
P6	886.07	752.01	1002.44	929.35	989.71	985.68	886.86	915.18	964.88	887.50
P7	1826.84	1906.38	1796.41	1614.83	1765.55	1770.15	1687.68	1854.18	1854.44	1960.58
P8	520.25	468.31	363.38	392.53	354.66	323.34	325.04	247.09	234.69	211.06
R-Squared	0.91	0.89	0.90	0.92	0.92	0.92	0.90	0.82	0.88	0.82

the network, but its SHAP values began to trend positively after 2014.

The importance of node attribute indicators, exemplified by 'PageRank', has been diminishing, contrasting sharply with the growing attractiveness of most POI indicators. This shift suggests a transition from region-based governance to function-oriented management.

Prior to 2020, the 'Enterprise' factor maintained a consistently positive relationship with metro ridership. However, this association reversed from 2020 onward, possibly reflecting the widespread shift to remote work during the COVID-19 pandemic. In contrast, the 'Residence' factor exerted only minimal influence before 2014, but subsequently grew increasingly polarized, ultimately manifesting a distinctly negative contribution. Consistent with common sense, as an indirect indicator of surrounding income levels, the 'Sales' factor largely suppressed passenger volumes, likely because more affluent populations have a wider range of transport alternatives at their disposal. Meanwhile, despite its relatively modest impact, the 'Shopping' factor sustained a long-term positive correlation with ridership throughout the observation period.

The aforementioned study reveals that the influence direction of NP indicators on station ridership predictions may undergo shifts over time. Additionally, certain indicators exhibit a diminishing overall impact on ridership, regardless of whether their influence is positive or negative. This trend is manifested by the gradual convergence of SHAP values' beeswarm plot scatter points from both sides towards the baseline, as observed with the 'Rent' indicator.

Due to the simultaneous presence of positive and negative influences, we determined the relative contribution of each indicator to ridership predictions by calculating the average of the absolute SHAP values corresponding to each feature. The quantified contributions of the indicators, based on their absolute values, are presented in Table III. Notably, the R^2 values of the predictive models across different years predominantly around 0.9, indicating excellent model fit and robustness.

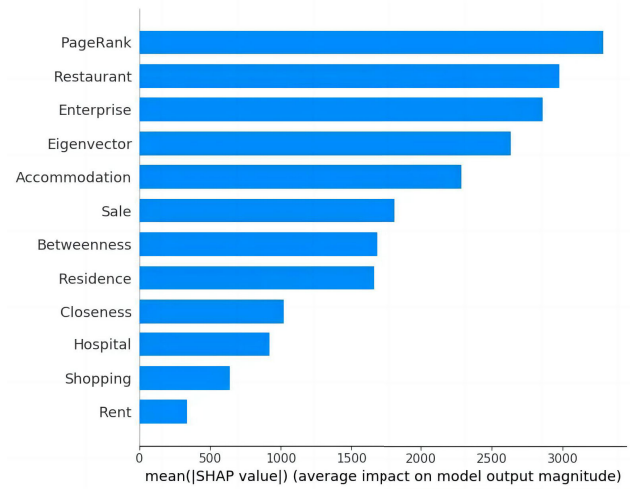


Fig. 6. Average absolute SHAP values' impact over a decade.

Fig. 6 illustrates the average absolute SHAP values of various features over the past decade. It is evident that the 'PageRank' indicator consistently occupies a prominent position in its contribution to the prediction outputs, irrespective of the direction of its influence. This underscores the effectiveness of node network location information in accurately forecasting ridership distribution.

Subsequently, the 'Restaurant' and 'Enterprise' index also play pivotal roles, with their contributions significantly surpassing those of other POI attributes. These findings highlight the strong association between the number of food and enterprise POIs in the vicinity of a station and its passenger ridership. Conversely, the contribution of the Rent indicator has a minimal impact on the prediction outcomes, rendering its influence on ridership almost negligible.

As illustrated in Fig. 7, the importance of contributing indicators remains dynamically variable over time. The assessment of each indicator's actual impact on the model's predictions, quantified by SHAP values, is continuously updated,

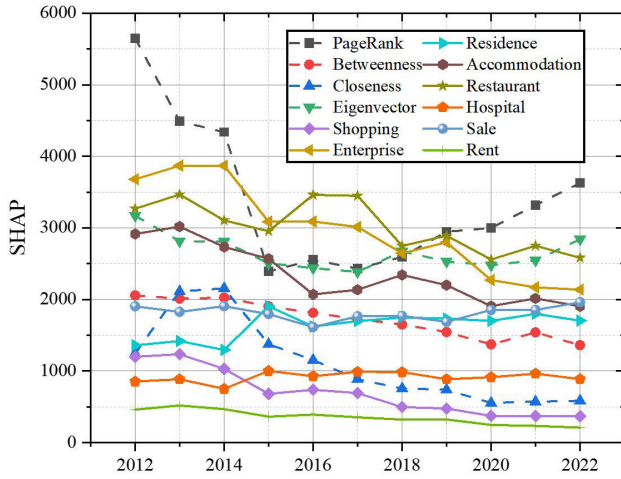


Fig. 7. Evolution of the SHAP values over a decade.

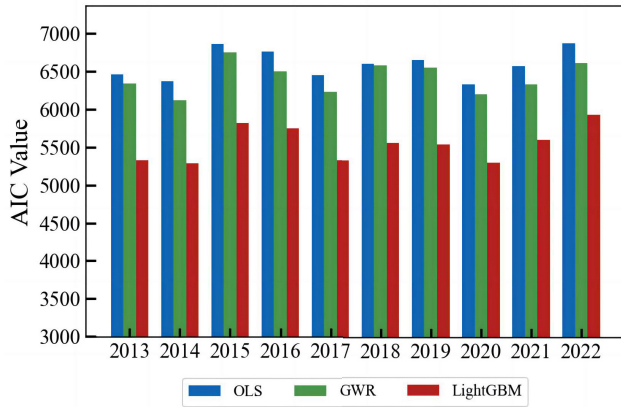


Fig. 8. AIC values among 2013-2022 when using OLS, GWR and LightGBM, respectively.

reflecting the evolving significance of different factors. Consequently, the model consistently maintains a high explanatory R^2 , ensuring robust and interpretable predictive outcomes. This dynamic adjustment underscores the adaptability of the machine learning model in capturing the shifting influences of various Node-Place indicators on metro ridership over the past decade.

C. Evaluation of Proposed Method

To appropriately evaluate the performance of proposed model, we employ the Akaike Information Criterion (based on AIC Value) [33], along with widely recognized metrics R^2 and adjusted R^2 . The AIC value is calculated as follows:

$$AIC = 2Q - 2LL(\theta) \quad (6)$$

where Q is the number of estimated parameters, and $LL(\theta)$ is the log-likelihood at the optimal parameter values θ . A lower AIC indicates a better trade-off between model fit and complexity, penalizing models with excessive parameters.

Leveraging Node-Place indicators, we apply the multiple regression method based on LightGBM to forecast the annual average daily entry passenger flow at Beijing metro stations. As shown in Fig. 8 and Fig. 9, the proposed method consistently outperforms both OLS and GWR in terms of AIC value, R^2 , and adjusted R^2 .

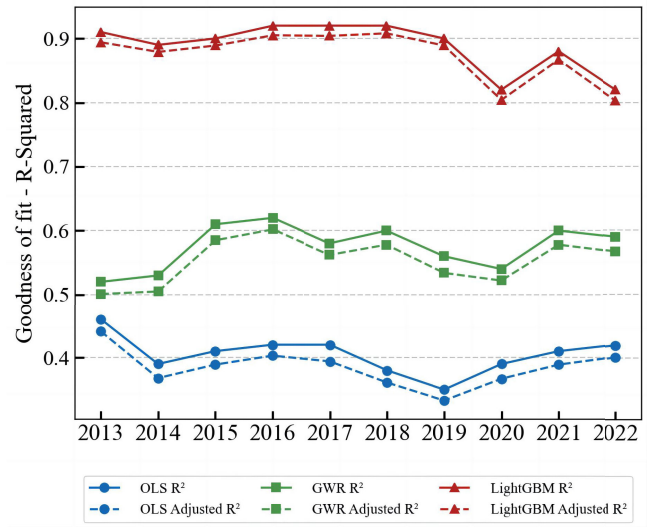


Fig. 9. Goodness of fit among 2013-2022.

TABLE IV
PARAMETERS FOR THE PROPOSED MODEL

Parameter	Description	Optimum
boosting_type	The type of boosting algorithm. (rf, dart, gbdt, et al)	gbdt
objective	The objective function determines the task type. (regression / classification)	regression
max_depth	The maximum depth of the decision tree. Limiting depth helps prevent overfitting, with a typical range of 3 to 10.	3
n_estimators	The number of weak learners, i.e., the number of trees in the ensemble. (Typically ranges from 100 to 1000)	100
learning_rate	The rate at which the algorithm moves in one step. (Typically ranges from 0.05 to 0.1)	0.1

Specifically, as illustrated in Fig. 8, the AIC values for OLS and GWR remain closely aligned across all examined years, with GWR consistently yielding marginally lower values than OLS. The LightGBM model significantly outperforms both, demonstrating AIC values that are consistently around 15% lower than those of OLS and GWR.

As depicted in Fig. 9, the LightGBM+SHAP machine learning framework introduced in our study not only maintains higher R^2 values but also exhibits a smaller gap between R^2 and adjusted R^2 compared to Ordinary Least Squares and Geographically Weighted Regression models [34]. Specifically, the results, as detailed in Table III, reveal an average R-squared value exceeding 0.88, significantly outperforming GWR's R^2 of 0.58 and OLS's R^2 of 0.42.

This consistently outstanding performance highlights LightGBM's ability to model ridership dynamics efficiently without overfitting.

TABLE IV presents the parameters set in proposed model, along with descriptions and optimal values. The boosting_type

TABLE V
MULTICOLLINEARITY DIAGNOSTICS OF VARIABLES USING VIF METHOD

VIF	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022
N1	3.12	2.72	2.44	1.26	1.19	1.25	1.21	1.25	1.35	1.43
N2	1.66	1.46	1.33	1.41	1.35	1.37	1.28	1.22	1.21	1.17
N3	2.68	3.60	3.45	2.14	1.81	1.90	2.30	2.31	1.76	1.91
N4	2.39	2.75	3.05	2.71	2.09	1.83	1.84	1.86	1.69	1.76
P1	4.12	3.69	3.00	1.37	1.71	1.60	1.40	1.35	1.23	1.16
P2	2.83	2.58	2.17	1.36	2.84	2.88	2.65	2.57	2.07	1.95
P3	2.43	2.71	2.48	1.77	1.88	1.73	2.22	2.13	3.12	2.64
P4	5.55	6.53	3.82	2.32	1.72	1.85	2.27	1.96	2.05	2.26
P5	6.40	6.50	4.35	2.78	3.29	3.39	2.91	2.91	3.31	2.89
P6	2.00	2.13	1.91	1.13	1.11	1.21	1.12	1.06	1.12	1.22
P7	3.72	2.46	2.10	1.73	1.61	1.69	1.71	1.67	2.48	2.37
P8	1.14	1.51	1.44	1.30	1.45	1.23	1.35	1.29	1.26	1.13

was set to gbdt, utilizing gradient boosting for stability, while the objective function was selected for regression, aligned with the task in this paper. Three other parameters were chosen to balance model complexity and efficiency. These parameters were determined through grid search to ensure optimal performance and computational efficiency.

In terms of computational efficiency, the proposed model and the benchmark models were implemented in the Python 3.9.13 environment on a laptop equipped with an Intel(R) Core(TM) i9-14900HX CPU and an NVIDIA GeForce RTX 4070 Laptop GPU.

Prior to the construction of analytical models (such as the Node-Place framework), as shown in TABLE V, a rigorous examination for multicollinearity among the variables is conducted by employing the Variation Inflation Factor (VIF) [35]. The results of this examination are reassuring, for all VIF values substantially below the threshold of 10, and the majority notably under 3. These findings substantiate the absence of multicollinearity concerns within the dataset, affirming the reliability of the variables for subsequent modeling. Furthermore, as illustrated in the evolutionary trajectories presented in Fig. 5 and Fig. 7, we conduct an additional analysis for the year 2012 to capture the transitional changes between 2012 and 2013. This supplementary analysis of 2012 is not included in the computation of the average absolute SHAP values depicted in Fig. 6.

V. DISCUSSION AND CONCLUSION

This study introduces an advanced machine learning framework based on SHAP that significantly surpasses traditional Ordinary Least Squares and Geographically Weighted Regression models in both regression performance and interpretability within the context of Beijing's metro network. By integrating diverse Node-Place indicators, the proposed model effectively captures the intricate relationships between metro network expansion and passenger ridership dynamics, offering a more nuanced analysis than previous methodologies. The performance of the SHAP framework underscore

its robustness and reliability, demonstrating its potential to enhance urban transportation and land use planning and inform policy development. Moreover, the application of SHAP values facilitates a deeper understanding of the specific contributions of various NP indicators, highlighting the pivotal role of node-centrality metrics such as PageRank and the significant impact of Points of Interest like restaurants and enterprises on ridership patterns.

Extending proposed methodology to other cities may further enhance its utility, providing new insights for transportation authorities seeking to improve the efficiency, sustainability, and economic impact of metro systems. In cities with limited transit infrastructure or varying demand patterns, the ability to accurately estimate ridership can lead to cost-effective expansion strategies, targeted investments in high-demand areas, and more sustainable transportation solutions.

Despite these advancements, the study acknowledges certain limitations. The reliance on available POI data may omit other influential factors such as real-time traffic conditions and seasonal variations, which could further refine ridership predictions. Additionally, while Beijing serves as a robust case study, the universality of our findings to other metropolitan areas requires further exploration.

Future research can benefit from incorporating more comprehensive map data like Areas of Interest (AOI) data, which is likely to better align with the evolving needs of related studies. Additionally, the impact of demographic factors and fare structures may prove to be significant, warranting further exploration. Another promising direction is the study of spatial features associated with the attraction side (destination of a trip) of metro ridership, which holds substantial value for understanding the complex dynamics of passenger flow and improving transit system planning.

REFERENCES

- [1] D. Sun, C. Zhang, M. Zhao, L. Zheng, and W. Liu, "Traffic congestion pattern detection using an improved mcmaster algorithm," in *Proc. 29th Chin. Control Decis. Conf. (CCDC)*, Chongqing, China, May 2017, pp. 2814–2819.

- [2] Q. Liu, E. Chung, and L. Zhai, "Fusing moving average model and stationary wavelet decomposition for automatic incident detection: Case study of Tokyo expressway," *J. Traffic Transp. Eng., English Ed.*, vol. 1, no. 6, pp. 404–414, 2014.
- [3] M. K. Hasan, "A framework for intelligent decision support system for traffic congestion management system," *Engineering*, vol. 2, no. 4, pp. 270–289, 2010.
- [4] K. Lu, B. Han, F. Lu, and Z. Wang, "Urban rail transit in China: Progress report and analysis (2008–2015)," *Urban Rail Transit*, vol. 2, nos. 3–4, pp. 93–105, Dec. 2016.
- [5] B. Han et al., "Statistical analysis of urban rail transit operation in the world in 2022: A review," *Urban Rapid Rail Transit*, vol. 36, no. 1, pp. 1–8, 2023.
- [6] D. Lin, J. D. Nelson, M. Beecroft, and J. Cui, "An overview of recent developments in China's metro systems," *Tunnelling Underground Space Technol.*, vol. 111, May 2021, Art. no. 103783.
- [7] China Urban Rail Transit Association. (Jan. 11, 2023). *Overview Mainland China's Urban Rail Transit Lines [EB/OL]*. Accessed: Jan. 3, 2023. [Online]. Available: <https://www.camet.org.cn/xyxw/11484>
- [8] A. Chakraborty and S. Mishra, "Land use and transit ridership connections: Implications for state-level planning agencies," *Land Use Policy*, vol. 30, no. 1, pp. 458–469, Jan. 2013.
- [9] X. Huang, Q. Liang, Z. Feng, and S. Chai, "A TOD planning model integrating transport and land use in urban rail transit station areas," *IEEE Access*, vol. 9, pp. 1103–1115, 2021.
- [10] Z.-J. Wang, F. Chen, B. Wang, and J.-L. Huang, "Passengers' response to transit fare change: An ex post appraisal using smart card data," *Transportation*, vol. 45, no. 5, pp. 1559–1578, Sep. 2018.
- [11] L. Bertolini, "Nodes and places: Complexities of railway station redevelopment," *Eur. Planning Stud.*, vol. 4, no. 3, pp. 331–345, 1996.
- [12] S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Proc. 31st Int. Conf. Neural Inf. Process. Syst. (NIPS'17)*. Red Hook, NY, USA: Curran Associates, pp. 4768–4777.
- [13] Y. Zhu, F. Chen, Z. Wang, and J. Deng, "Spatio-temporal analysis of rail station ridership determinants in the built environment," *Transportation*, vol. 46, no. 6, pp. 2269–2289, Dec. 2019.
- [14] S. Liu, E. Yao, and B. Li, "Exploring urban rail transit station-level ridership growth with network expansion," *Transp. Res. D, Transp. Environ.*, vol. 73, pp. 391–402, Aug. 2019.
- [15] L. Bertolini, *Station Areas as Nodes and Places in Urban Networks: An Analytical Tool and Alternative Development Strategies*. Physica-Verlag HD, 2008, pp. 35–57, doi: [10.1007/978-3-7908-1972-4_3](https://doi.org/10.1007/978-3-7908-1972-4_3).
- [16] L. Bertolini, "Spatial development patterns and public transport: The application of an analytical model in The Netherlands," *Planning Pract. Res.*, vol. 14, no. 2, pp. 199–210, May 1999.
- [17] Y. Zhao, S. Hu, and M. Zhang, "Evaluating equitable transit-oriented development (TOD) via the node-place-people model," *Transp. Res. A, Policy Pract.*, vol. 185, Jul. 2024, Art. no. 104116.
- [18] M. Zhou, J. Zhou, J. Zhou, S. Lei, and Z. Zhao, "Introducing social contacts into the node-place model: A case study of Hong Kong," *J. Transp. Geography*, vol. 107, Feb. 2023, Art. no. 103532.
- [19] G. Lyu, L. Bertolini, and K. Pfeffer, "Developing a TOD typology for Beijing metro station areas," *J. Transp. Geography*, vol. 55, pp. 40–50, Jul. 2016.
- [20] A. Oseni et al., "An explainable deep learning framework for resilient intrusion detection in IoT-enabled transportation networks," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 1, pp. 1000–1014, Jan. 2023, doi: [10.1109/TITS.2022.3188671](https://doi.org/10.1109/TITS.2022.3188671).
- [21] E. Štrumbelj and I. Kononenko, "Explaining prediction models and individual predictions with feature contributions," *Knowl. Inf. Syst.*, vol. 41, no. 3, pp. 647–665, Dec. 2014.
- [22] M. Li, H. Sun, Y. Huang, and H. Chen, "SVCE: Shapley value guided counterfactual explanation for machine learning-based autonomous driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 10, pp. 14905–14916, Oct. 2024.
- [23] C. Fu, Z. Huang, B. Scheuer, J. Lin, and Y. Zhang, "Integration of dockless bike-sharing and metro: Prediction and explanation at origin-destination level," *Sustain. Cities Soc.*, vol. 99, Dec. 2023, Art. no. 104906.
- [24] H. Li, Y. Jin, and G. Ren, "Interpretable prediction of pedestrian crossing intention: Fusion of human skeletal information in natural driving scenarios," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 11, pp. 18153–18170, Nov. 2024.
- [25] M. Li, M.-P. Kwan, W. Hu, R. Li, and J. Wang, "Examining the effects of station-level factors on metro ridership using multiscale geographically weighted regression," *J. Transp. Geography*, vol. 113, Dec. 2023, Art. no. 103720.
- [26] H. Shi, Y. Li, Z. Jiang, and J. Yan, "Comprehensive evaluation of power quality for microgrid based on CRITIC method," in *Proc. IEEE 9th Int. Power Electron. Motion Control Conf. (IPEMC-ECCE Asia)*, Nanjing, China, Nov. 2020, pp. 1667–1669.
- [27] H. Yang, D. Du, J. Wang, X. Wang, and F. Zhang, "Reshaping China's urban networks and their determinants: High-speed rail vs. air networks," *Transp. Policy*, vol. 143, pp. 83–92, Nov. 2023.
- [28] Z. Li, J. Tang, C. Zhao, and F. Gao, "Improved centrality measure based on the adapted PageRank algorithm for urban transportation multiplex networks," *Chaos, Solitons Fractals*, vol. 167, Feb. 2023, Art. no. 112998.
- [29] E. Chen, Y. Liu, and M. Yang, "Revealing senior mobility patterns and activities in urban transit systems," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 10, pp. 11424–11437, Oct. 2023.
- [30] W. Luo, J. Liu, and X. Xu, "Examining the relationship between built environment and metro ridership at zone-to-zone level," in *Proc. IEEE 26th Int. Conf. Intell. Transp. Syst. (ITSC)*, Bilbao, Spain, Sep. 2023, pp. 2269–2274.
- [31] G. Ke et al., "LightGBM: A highly efficient gradient boosting decision tree," in *Proc. 31st Int. Conf. Neural Inf. Process. Syst. (NIPS'17)*. Red Hook, NY, USA: Curran Associates, pp. 3149–3157.
- [32] Y. Jing, H. Hu, S. Guo, X. Wang, and F. Chen, "Short-term prediction of urban rail transit passenger flow in external passenger transport hub based on LSTM-LGB-DRS," *IEEE Trans. Intell. Transp. Syst.*, vol. 22, no. 7, pp. 4611–4621, Jul. 2021.
- [33] E. Chen, Z. Ye, C. Wang, and M. Xu, "Subway passenger flow prediction for special events using smart card data," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 3, pp. 1109–1120, Mar. 2020, doi: [10.1109/TITS.2019.2902405](https://doi.org/10.1109/TITS.2019.2902405).
- [34] A. E. Iyanda and T. Osayomi, "Is there a relationship between economic indicators and road fatalities in texas? A multiscale geographically weighted regression analysis," *GeoJournal*, vol. 86, no. 6, pp. 2787–2807, Dec. 2021.
- [35] R. L. Mason, R. F. Gunst, and J. L. Hess, *Statistical Design and Analysis of Experiments: With Applications To Engineering and Science*. Hoboken, NJ, USA: Wiley, 2003.



Yizhe Wang received the B.E. degree from the Jeme Tienyow Honors College, Beijing Jiaotong University, in 2024. He is currently pursuing the Ph.D. degree with the School of Civil Engineering, Beijing Jiaotong University.

His main research interests include urban rail transit based on multisource big data and artificial intelligence.



Zijia Wang received the Ph.D. degree from Beijing Jiaotong University in 2013. He is currently a Professor with the School of Civil Engineering, Beijing Jiaotong University. His research interests include multidimensional data mining in rail transit big data and the application of virtual reality (VR) technology in rail transit.



Fanxi Zhao is currently pursuing the bachelor's degree in actuarial science and statistics with the University of Waterloo, Waterloo, ON, Canada. Her main research interests include the application of machine learning in the actuarial field, including its role in financial modeling and actuarial processes.